



cimt

let your work flow

**COMPETENCE FOR BI, BIG DATA AND ESB PROJECTS**

How to steer and monitor jobs. Pattern and Best Practices.

# Agenda

- 1 Different Kind of Data Integration Jobs
- 2 Requirements
- 3 Overview about Components and Data Model
- 4 Using Components for Simple Jobs
- 5 Using Components for Incremental Loads
- 6 Using Components for Partinioned Data Loads
- 7 Monitoring Reports
- 8 Logging
- 9 Outlook

# Different Kind of Data Integration Jobs

## » Simple Jobs:

- Does not need input parameters to know which data should be processed
- Does not start on a particular piece of data
- Does not need to partition the data to process
- Example: Any job

## » Incremental Loads

- Proceed data which follows the data previously proceed
- Depends and a sequential criteria to know where to start
- Example: Growing source tables

## » Partitioned Data Loads

- Data has a defined start and end
- Data items to process separately
- Data are pre-packaged / partitioned
- Example: File processing / Calendar based data processing

# Agenda

- 1 Different Kind of Data Integration Jobs
- 2 Requirements
- 3 Overview about Components and Data Model
- 4 Using Components for Simple Jobs
- 5 Using Components for Incremental Loads
- 6 Using Components for Partinioned Data Loads
- 7 Monitoring Reports
- 8 Logging
- 9 Outlook

# Requirements for Data Integration Jobs

- » Restart ability
  - Jobs should detect it self what is the work to do
  - Jobs should finish its work completely or nothing
  - It must be allowed to start a job without problems
- » Limited Resource Allocation
  - Jobs should not need more the 1GB RAM
  - Jobs have to process huge amount of data in a streaming like way
  - Jobs should use a closed data range to process
- » Projectable processing
  - It is by far more important to have a system in which the administrator can monitor the progress than a system which is some percentage faster but does not allow the recognise any progress.
  - It is also a design goal to spread the workload over a longer period of time instead of having a huge peak which makes the system a irresponsible for a unknown period of time.

# Agenda

- 1 Different Kind of Data Integration Jobs
- 2 Requirements
- 3 Overview about Components and Data Model
- 4 Using Components for Simple Jobs
- 5 Using Components for Incremental Loads
- 6 Using Components for Partinioned Data Loads
- 7 Monitoring Reports
- 8 Logging
- 9 Outlook

# Talend Components and Data Model

- » The system consists of 4 Talend components:
  - tJobInstanceStart:
    - Creates an entry in the JOB\_INSTANCE\_STATUS table when the job starts.
    - Initialises the logging framework
    - Provides information about previous job runs
  - tJobInstanceEnd:
    - Collects measurements
    - Collects status information about the job run
    - Updates the entry in the JOB\_INSTANCE\_STATUS table
  - tJobDataRangeScanner:
    - Scans a flow for min/max values to set them as time range start/end and value range start/end.
  - tJobInstanceLifeCheck
    - Checks if some open entries in the JOB\_INSTANCE\_STATUS table are not alive anymore and close them up.
- » All components are described in detail in the PDF documentation

# Talend Components and Data Model

- » The data model consists of 4 core tables and one optional table:
  - JOB\_INSTANCE\_STATUS
    - Job meta data
    - Timestamps
    - Counters
    - OS information
    - Return code and messages
  - JOB\_INSTANCES\_LOGS
    - Log messages
  - JOB\_INSTANCE\_COUNTERS
    - Detail counters
  - JOB\_INSTANCE\_CONTEXT
    - Context variables at the start and at the end of the job run.
  - JOB\_CALENDAR
    - A full fledged calendar with financial date ranges



## Architecture hints

- » Take care of the kind of generation of the JOB\_INSTANCE\_ID
- » Add the column JOB\_INSTANCE\_ID (case sensitive is not an issue) to the actual data tables to provide data lineage.
- » Sometimes it is supposed to have 2 different columns job\_instance\_id\_insert and job\_instance\_id\_update if you want to track which job has created and updated the data.
- » Use the field option of the output components to send the job\_instance\_id according to the statement type.
- » Decide if the table names fits to the customers naming rules. If not set alternative table names.
- » Decide if the amount of job runs are very high and use more than one schema + tables for the monitoring. E.g. use a dedicated schema + tables for the staging jobs and another schema + tables for the core and mart jobs.
- » Another approach is to establish different schemas for different (not interoperating) parts of the system.
- » Probably it is supposed to have partitioned table for JOB\_INSTANCE\_STATUS (e.g. with the month as partitions or use the project as partitioning criteria)

# Agenda

- 1 Different Kind of Data Integration Jobs
- 2 Requirements
- 3 Overview about Components and Data Model
- 4 Using Components for Simple Jobs
- 5 Using Components for Incremental Loads
- 6 Using Components for Partinioned Data Loads
- 7 Monitoring Reports
- 8 Logging
- 9 Outlook

# Using the framework for a simple job

- » It provides runtime information
- » Allows to report the job runs in a unified way

The screenshot displays a job designer interface with two workflow diagrams. The top diagram shows a 'Prejob' section with components tPrejob\_2, DWH\_META, and tJobInstanceStart\_1, followed by a 'Postjob' section with tPostjob\_2 and tJobInstanceEnd\_1. The bottom diagram shows a sequence of DWH\_MOBILE, DWH\_MOBILE\_MIRROR, and tMysqlTableTransfer\_1. Below the diagrams is a configuration panel for the 'tJobInstanceEnd\_1' component. The 'Input counters' section is highlighted with an orange border and contains the following table:

Input counter (use NB_LINE etc.)	Add (unchecked for minus)	Name
{{(Integer)globalMap.get("tMysqlTableTransfer_1_NB_LINE")}}	<input checked="" type="checkbox"/>	

The 'Output counters' section contains the following table:

Output counter (use NB_LINE etc.)	Add (unchecked for minus)	Name
{{(Integer)globalMap.get("tMysqlTableTransfer_1_NB_INSERT")}}	<input checked="" type="checkbox"/>	

Counters:  
Use the return values  
from the components.  
Hint: Use CTRL+Space

# Agenda

- 1 Different Kind of Data Integration Jobs
- 2 Requirements
- 3 Overview about Components and Data Model
- 4 Using Components for Simple Jobs
- 5 Using Components for Incremental Loads
- 6 Using Components for Partinioned Data Loads
- 7 Monitoring Reports
- 8 Logging
- 9 Outlook

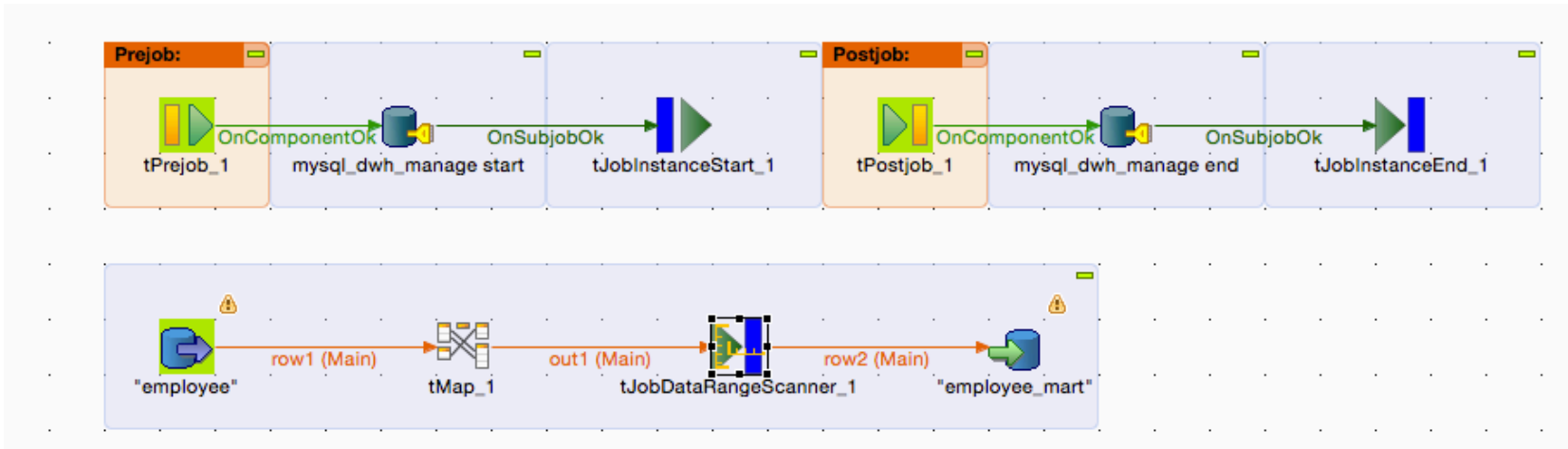
# Using the framework for incremental steering

## Source column based

- » Use the data ranges from the previous run to select the input data for the current run
- » Decide which kind of value do you can use to identify the new / changed records in the source
  - Time based ranges (Date type)
  - Value based ranges (Long or String typed)
- » Detect the range within the job run with `tJobDataRangeScanner` (this avoids to measure it with a probably expensive select in the source or target table).
- » Especially for time ranges it is very important NOT to rely on the local job run time. The time on different servers potentially differs.
- » We have to measure the last time stamp within the data because after the job has been finished, it is possible to get data with an timestamp BEFORE the timestamp of the job end – e.g. because of delayed mirroring mechanism.

# Using the framework for incremental steering – job design

## Source column based



It is important to measure the data ranges near the target to be sure these Records are persisted. A good idea is also using a transaction.

Information about previous job runs

Return last instance results  Last successful  Last must have data inserted or deleted  For the current work item

Collect Job Instance IDs running after previous run

OK Result Codes (Comma separated list) "0,100"

Replacement for prev. job-instance-id if job runs first time 01

Replacement for prev. job-start-date if job runs first time TalendDate.parseDate("yyyy-MM-dd", "1970-01-01")

Replacement for prev. time-range-end if job runs first time TalendDate.parseDate("yyyy-MM-dd", "1970-01-01")

Replacement for prev. value-range-end if job runs first time "0"

Replacement for prev. result-item if job runs first time "0"

Necessary settings of the `tJobInstanceStart` component to get the information about the last job run.

# Using the framework for incremental steering – job design

## Source column based

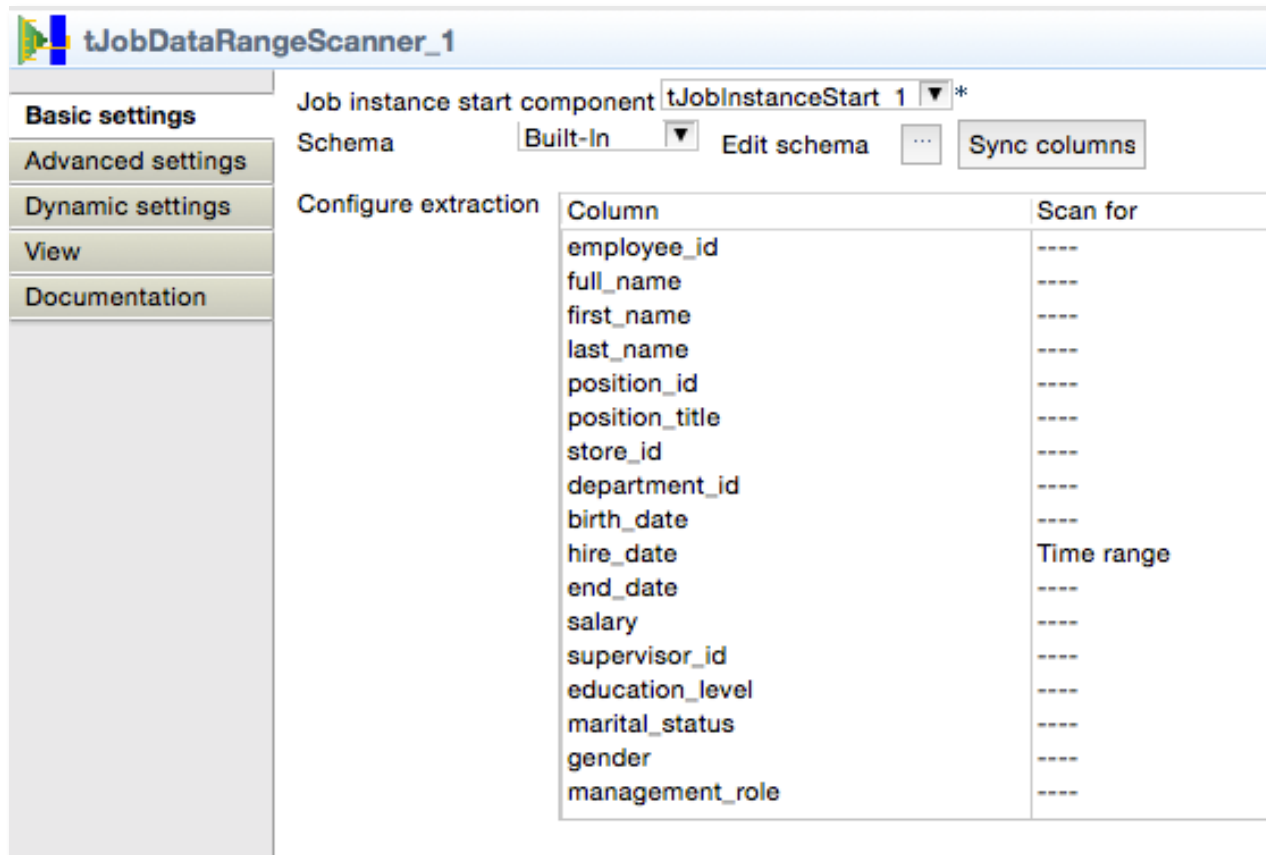
Table Name	"employee"
Query Type	Built-In <input type="button" value="Guess Query"/> <input type="button" value="Guess schema"/>
Query	<pre>"SELECT \'+context.postgres_foodmart_Database+"\.'public\.'employee\.'employee_id\', \'+context.postgres_foodmart_Database+"\.'public\.'employee\.'full_name\', \'+context.postgres_foodmart_Database+"\.'public\.'employee\.'first_name\', \'+context.postgres_foodmart_Database+"\.'public\.'employee\.'last_name\', \'+context.postgres_foodmart_Database+"\.'public\.'employee\.'position_id\', \'+context.postgres_foodmart_Database+"\.'public\.'employee\.'position_title\', \'+context.postgres_foodmart_Database+"\.'public\.'employee\.'store_id\', \'+context.postgres_foodmart_Database+"\.'public\.'employee\.'department_id\', \'+context.postgres_foodmart_Database+"\.'public\.'employee\.'birth_date\', \'+context.postgres_foodmart_Database+"\.'public\.'employee\.'hire_date\', \'+context.postgres_foodmart_Database+"\.'public\.'employee\.'end_date\', \'+context.postgres_foodmart_Database+"\.'public\.'employee\.'salary\', \'+context.postgres_foodmart_Database+"\.'public\.'employee\.'supervisor_id\', \'+context.postgres_foodmart_Database+"\.'public\.'employee\.'education_level\', \'+context.postgres_foodmart_Database+"\.'public\.'employee\.'marital_status\', \'+context.postgres_foodmart_Database+"\.'public\.'employee\.'gender\', \'+context.postgres_foodmart_Database+"\.'public\.'employee\.'management_role\' FROM \'+context.postgres_foodmart_Database+"\.'public\.'employee\' where hire_date &gt; "' + TalendDate.formatDate("yyyy-MM-dd", ((java.util.Date)globalMap.get("tJobInstanceStart_1_PREV_TIME_RANGE_END"))) + "'"</pre>

Settings of the source database input component.

Here we use the time range end value of the previous run to select the data newer as we already got in the last run.

# Using the framework for incremental steering – job design

## Source column based



The component tJobDataRangeScanner does not have any impact on its flow. It simply scans columns on the min/max value or timestamps and sets these values as time or value range min/max values.



# Using the framework for incremental steering – job design

## Job\_instance\_id based

tJobInstanceStart settings to return the list of job instance ids of all job runs after the previous run of this job.

Information about previous job runs

Return last instance results  Last successful  Last must have data inserted or deleted  For the current work item

Collect Job Instance IDs running after previous run  Only successful  Only with data

Source job names	Job name
	"stage_import_job"

If the source and target table does have job\_instance\_id columns, these column can be used to detect changed records.

The new or changed source records can be selected with:

```
where job_instance_id in (" +  
((String)globalMap.get("tJobInstanceStart_1_SOURCE_JOB_INSTANCE_ID_LIST")) + ")"
```

With this approach it is possible to build a fully generic way to detect the necessary records to process in data ware house systems.

# Agenda

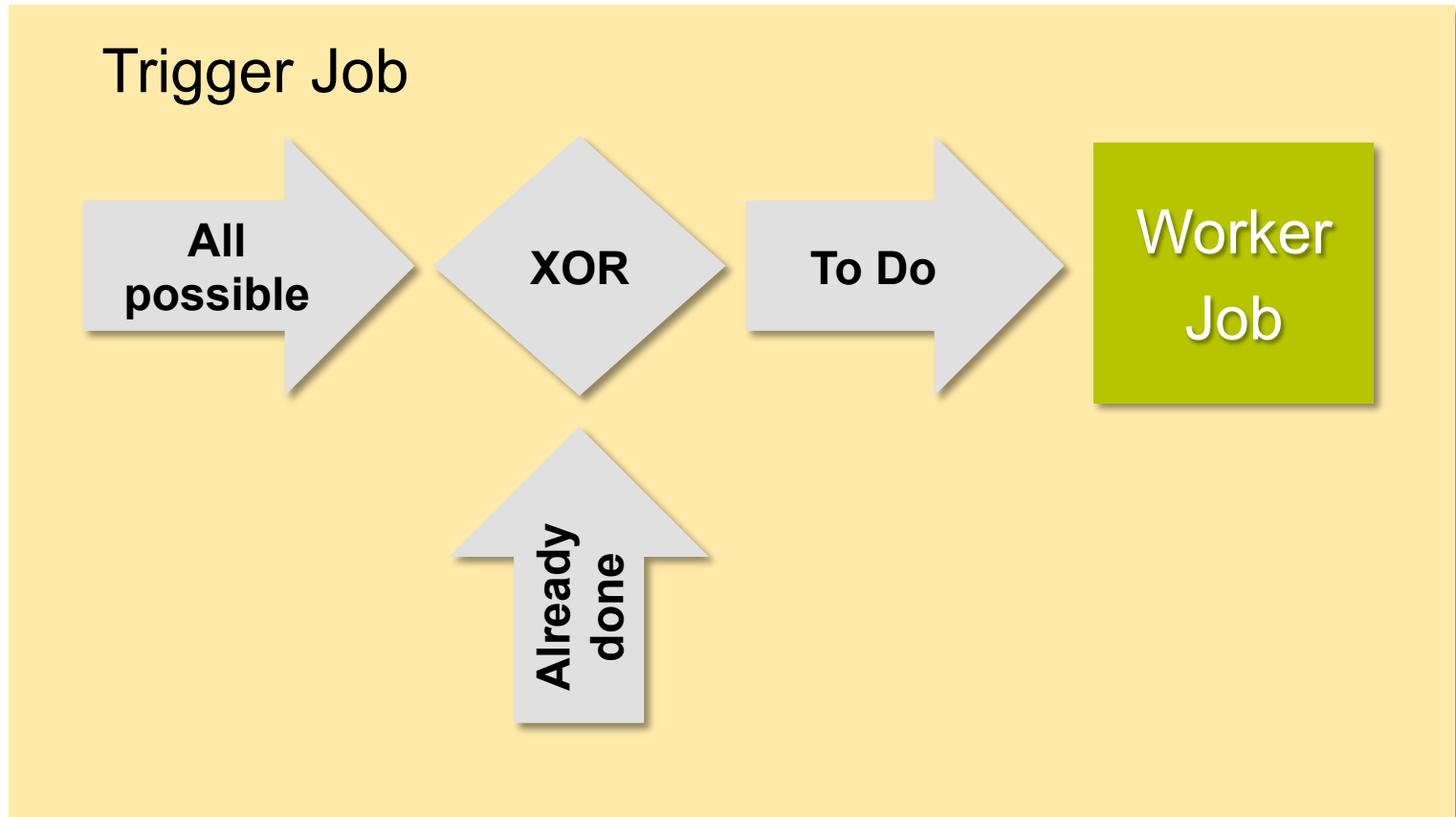
- 1 Different Kind of Data Integration Jobs
- 2 Requirements
- 3 Overview about Components and Data Model
- 4 Using Components for Simple Jobs
- 5 Using Components for Incremental Loads
- 6 Using Components for Partinioned Data Loads
- 7 Monitoring Reports
- 8 Logging
- 9 Outlook

## Using the framework for partitioned data loads

- » The source data have to have an dedicated realm
- » Examples:
  - Data are content of an file
  - Data are related to a date or any possible incremental value like countries, distinct ids.
- » On data element (file, date etc.) is called a work item
- » The target data should not rely on the previous loaded data.
- » The pattern needs to know all possible work items and the work items already successful processed.
- » The amount of possible work items can be become a huge number and should therefore reduced to a reasonable amount like the dates of the last year.
- » This pattern consists of 2 jobs: a trigger job which selects the work items to do and a worker job which process the data of one work item.
- » Pattern allows the parallel processing of work items as long as they can processed separately

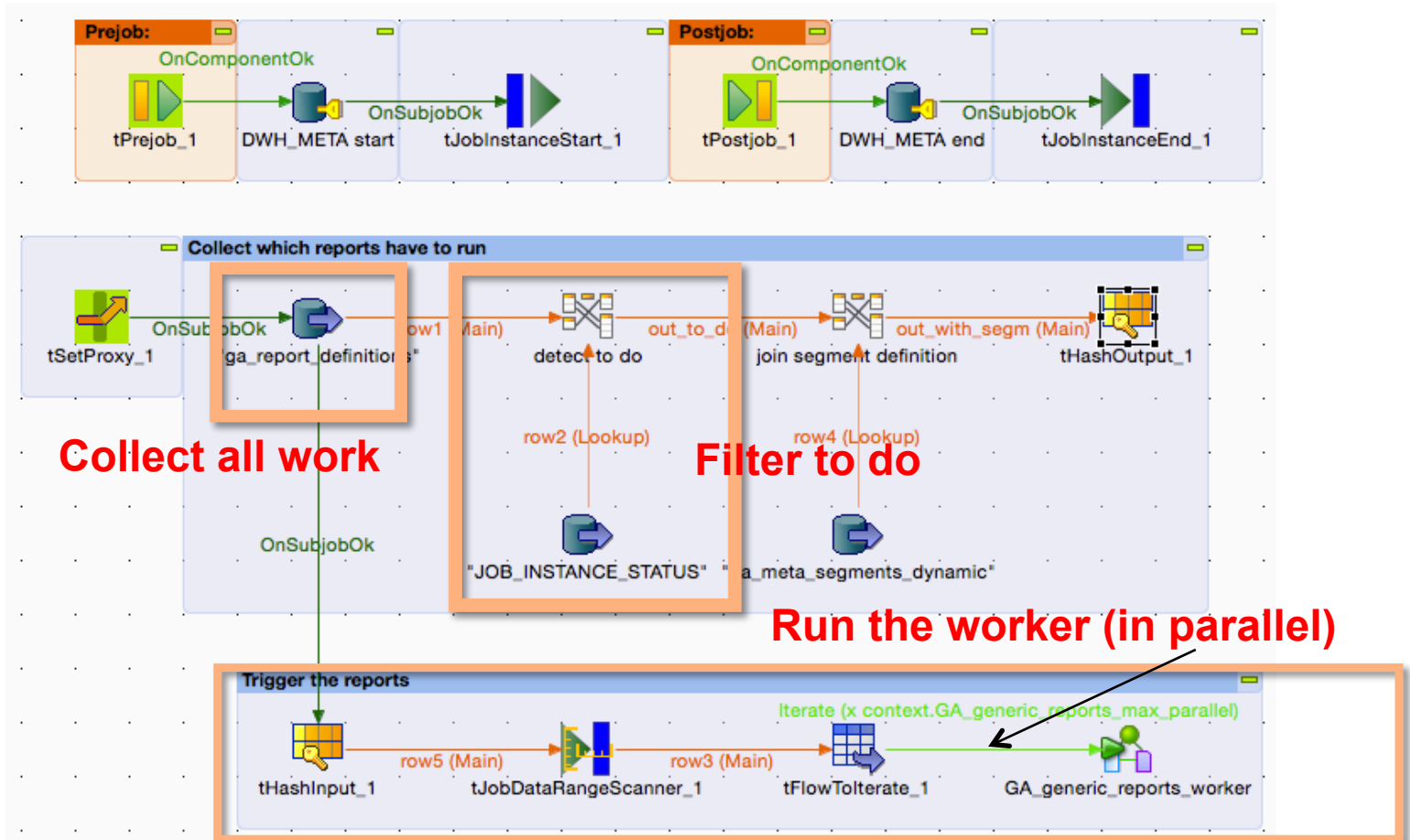
## Using the framework for partitioned data loads

- » Trigger and worker job pattern.



# Using the framework for partitioned data loads

- » Select the total amount of work items
- » Select the amount of done work items
- » Build the difference and start the worker job(s - in case of parallel runs)



# Using the framework for partitioned data loads

## Select the done work items

"JOB\_INSTANCE\_STATUS"(tMysqlInput\_2)

Basic settings

Advanced settings

Dynamic settings

View

Documentation

Property Type: Repository DB (MYSQL):DWH\_META

DB Version: Mysql 5

Use an existing connection

Host: context.DWH\_META\_Server \* Port: context.DWH\_META\_Port

Username: context.DWH\_META\_Login

Schema: Built-In Edit schema

Table Name: "JOB\_INSTANCE\_STATUS"

Query Type: Built-In Guess Query Guess schema

Query: "select work\_item from JOB\_INSTANCE\_STATUS where job\_name=" + jobName + "\_worker" and return\_code=0 order by work\_item desc"

If the worker job is named like trigger job + `_worker` it is very easy to select the done work items.

# Using the framework for partitioned data loads

## Select the work items to do

The screenshot displays the Talend Open Studio for Data Integration interface. The main workspace shows a tMap configuration with two rows: 'row1' and 'row2'. 'row1' is a data source with columns: report\_id, profile\_id, metrics, filters, dimensions, segment, and at\_date. 'row2' is a join component with the following properties:

Property	Value
Lookup Model	Load once
Match Model	Unique match
Join Model	Inner Join
Store temp data	false

The 'Expr. key' section shows the expression: `row1.report_id + "#" + TalendDate.formatDate("yyyyMMdd", row1.at_date)` mapped to the column `work_item`.

On the right, the 'out\_to\_do' table is visible, showing the 'Catch lookup inner join reject' property set to 'true'. Below it, the 'Expression' and 'Column' mapping is shown:

Expression	Column
row1.report_id	report_id
row1.profile_id	profile_id
row1.metrics	metrics
row1.filters	filters
row1.dimensions	dimensions
row1.segment	segment
row1.at_date	report_date

The 'Schema editor' at the bottom shows the 'row1' and 'out\_to\_do' schemas. The 'row1' schema has columns: report\_id (Key, int, Nullab, Date Pattern, Length 20, Precision 0), profile\_id (long, Nullab, Length 128, Precision 0), and metrics (String, Nullab, Length 128, Precision 0). The 'out\_to\_do' schema has columns: report\_id (Key, int, Nullab, Date Pattern, Length 20, Precision 0), profile\_id (long, Nullab, Length 128, Precision 0), and metrics (String, Nullab, Length 128, Precision 0).

The work items (total and done) will be joined and the not matching items are the work items to do. Please notice the inner join and the inner join reject option!

# Agenda

- 1 Different Kind of Data Integration Jobs
- 2 Requirements
- 3 Overview about Components and Data Model
- 4 Using Components for Simple Jobs
- 5 Using Components for Incremental Loads
- 6 Using Components for Partinioned Data Loads
- 7 Monitoring Reports
- 8 Logging
- 9 Outlook



# Monitoring Reports

Choose the parameters:

- Start and end date
- Filter job names Select the jobs to report (nothing select means all)
- Restrict to failed jobs
- Restrict to running jobs

The screenshot shows a dialog box titled "Eingabesteuerelemente" (Input Control Elements). It contains the following fields and options:

- \* Job Started Range Begin:** A date input field containing "23.06.2015".
- Job Started Range End:** An empty date input field.
- Job Filter:** A text input field containing "%".
- Jobs:** A list of job names with checkboxes next to them:
  - AI\_OEM\_Import\_master\_data
  - catch\_ddl\_data
  - cl7\_ACCOUNT\_export
  - cl7\_ACCOUNT\_STATUS\_HISTORY\_export
  - cl7\_all\_export\_new
  - cl7\_all\_export\_wrapper
  - cl7\_CAMPAIGN\_export
  - cl7\_CASE\_COMMENT\_export
  - cl7\_CASE\_export
  - cl7\_CASE\_HISTORY\_export
  - cl7\_CASE\_SOLUTION\_export
  - cl7\_CONTACT\_export
- Alle Keine Umkehren:** A label indicating that all jobs are selected by default.
- Show Only Errors:**  checkbox.
- Is Running:**  checkbox.

At the bottom of the dialog, there are four buttons: "Anwenden" (Apply), "OK", "Zurücksetzen" (Reset), and "Abbrechen" (Cancel).

# Monitoring Reports – Job List

Report list of selected jobs in the selected date range with the options.

ID	Root ID	Name	Started at	Ended at	Duration in s	Work Item	Time Range Start	Time Range End	Input	Output	Return Code	Message
1571488		f_daily_demand	2015-06-23 08:45:53	2015-06-23 08:47:26	93				0	0	0	
1571487	1571484	ci7_dwh_load_mirror_table	2015-06-23 08:45:14	2015-06-23 08:46:52	98	custom_campaignmember__c			0	0	0	
1571486	1571484	ci7_dwh_load_table	2015-06-23 08:45:13	2015-06-23 08:45:13	0	custom_campaignmember__c/mt/bimig			0	0	0	
1571485	1571484	ci7_dwh_expand_online_files	2015-06-23 08:45:12	2015-06-23 08:45:13	1	mt/bimigration/CRM/salesforce/import/2			0	0	0	
1571484		ci7_dwh_import_all_wrapper	2015-06-23 08:45:12	2015-06-23 08:59:20	848				0	0	0	
1571483	1571482	GA_generic_check_fact_table_columns	2015-06-23 08:44:12	2015-06-23 08:44:12	0				5	0	0	
1571482		GA_generic_facts_core	2015-06-23 08:44:12									
1571481	1570941	GA_generic_reports_worker	2015-06-23 08:43:43	2015-06-23 08:43:50	7	522#20150622	2015-06-22 00:00:00	2015-06-22 00:00:00	3	3	0	
1571480	1570941	GA_generic_reports_worker	2015-06-23 08:43:38	2015-06-23 08:43:52	14	520#20150622	2015-06-22 00:00:00	2015-06-22 00:00:00	3	3	0	
1571479	1570941	GA_generic_reports_worker	2015-06-23 08:43:41	2015-06-23 08:43:52	11	521#20150622	2015-06-22 00:00:00	2015-06-22 00:00:00	3	3	0	
1571478	1570941	GA_generic_reports_worker	2015-06-23 08:43:36	2015-06-23 08:43:52	16	519#20150622	2015-06-22 00:00:00	2015-06-22 00:00:00	3	3	0	
1571477	1570941	GA_generic_reports_worker	2015-06-23 08:43:34	2015-06-23 08:43:52	18	518#20150622	2015-06-22 00:00:00	2015-06-22 00:00:00	3	3	0	
1571476	1570941	GA_generic_reports_worker	2015-06-23 08:43:28	2015-06-23 08:43:36	8	517#20150622	2015-06-22 00:00:00	2015-06-22 00:00:00	3	3	0	
1571475	1570941	GA_generic_reports_worker	2015-06-23 08:43:26	2015-06-23 08:43:36	10	516#20150622	2015-06-22 00:00:00	2015-06-22 00:00:00	3	3	0	
1571474	1570941	GA_generic_reports_worker	2015-06-23 08:43:22	2015-06-23 08:43:36	14	514#20150622	2015-06-22 00:00:00	2015-06-22 00:00:00	3	3	0	
1571473	1570941	GA_generic_reports_worker	2015-06-23 08:43:19	2015-06-23 08:43:32	13	513#20150622	2015-06-22 00:00:00	2015-06-22 00:00:00	3	3	0	
1571472	1570941	GA_generic_reports_worker	2015-06-23 08:43:24	2015-06-23 08:43:34	10	515#20150622	2015-06-22 00:00:00	2015-06-22 00:00:00	3	3	0	
1571471	1570941	GA_generic_reports_worker	2015-06-23 08:43:16	2015-06-23 08:43:26	10	512#20150622	2015-06-22 00:00:00	2015-06-22 00:00:00	3	3	0	
1571470	1570941	GA_generic_reports_worker	2015-06-23 08:43:12	2015-06-23 08:43:19	7	511#20150622	2015-06-22 00:00:00	2015-06-22 00:00:00	3	3	0	
1571469	1570941	GA_generic_reports_worker	2015-06-23 08:43:07	2015-06-23 08:43:19	12	509#20150622	2015-06-22 00:00:00	2015-06-22 00:00:00	3	3	0	
1571468	1570941	GA_generic_reports_worker	2015-06-23 08:43:10	2015-06-23 08:43:19	9	510#20150622	2015-06-22 00:00:00	2015-06-22 00:00:00	3	3	0	
1571467	1570941	GA_generic_reports_worker	2015-06-23 08:42:58	2015-06-23 08:43:17	19	508#20150622	2015-06-22 00:00:00	2015-06-22 00:00:00	3	3	0	

Yellow highlighted jobs are running.  
Red highlighted jobs have been failed.

# Monitoring Reports – Job Details

Details are shown after follow the link on the Job Instance ID

The screenshot shows a web browser window with the following details for Job Instance ID 1567689:

- Job Instance ID: [1567689](#) Talend PID: nKreKr
- Technical Name: **ai\_smar\_fetch\_ads\_image\_url\_one\_customer**
- Display Name:
- Started At: 22.06.2015 10:14:43 Ended At: 22.06.2015 10:14:47 Duration: 4 s
- Root Job Instance: [1567547](#) Talend Root PID: 20150622101208\_4v0UB
- Job Info: AI\_CORE/ai\_smar\_fetch\_ads\_image\_url\_one\_customer-1.2:Prod
- Return Code: 4
- Job Server: talend38-8 OS PID: 28558
- OS User: tisadmin
- Work Item: 703114#15627
- Time Range Start: Time Range End:
- Value Range Start: Value Range End:
- Job Result:
- Count Input: 0 Count Output: 1 Count Updated: 0
- Count Rejects: 0 Count Deleted: 0
- Return Message: tDie\_1:Job failed. Last ad-id: 208578875  
tMysqlInput\_1:Communications link failure
- [View Log Table](#) The last packet sent successfully to the server was 0 milliseconds ago. The driver has not received any packets from the server.

# Agenda

- 1 Different Kind of Data Integration Jobs
- 2 Requirements
- 3 Overview about Components and Data Model
- 4 Using Components for Simple Jobs
- 5 Using Components for Incremental Loads
- 6 Using Components for Partinioned Data Loads
- 7 Monitoring Reports
- 8 Logging
- 9 Outlook

# Logging

- » Every job has its own logger
- » Every job instance has its own appender
  - Allows one file per job run
  - Allows distinct path per environment
  - Layout pattern can be vary from job to job if necessary
- » Appender are flushed and closed at the end of the job instance
  - Appender to database (JOB\_INSTANCE\_LOGS table)
  - Appender to files
    - Pattern layout can contain context variables
    - Log file path can be set based on context variables
- » Custom log4j.xml configuration can be used
- » See the component documentation for configuration details

# Agenda

- 1 Different Kind of Data Integration Jobs
- 2 Requirements
- 3 Overview about Components and Data Model
- 4 Using Components for Simple Jobs
- 5 Using Components for Incremental Loads
- 6 Using Components for Partinioned Data Loads
- 7 Monitoring Reports
- 8 Logging
- 9 Outlook

# Outlook

- » Additional table for work items within a job
  - Solves the problem for very frequent jobs processing a high number of work items with the fast growing values space of `job_instance_ids`
- » Web Service to concentrate the database actions.
  - Solves the problem of huge amount parallel database sessions and allows caching of pre selected job instance ids
- » Reporting will be enhanced with gant diagrams to display the scheduling
- » Dump of context variables (without passwords)
- » Make the Job Logger available for code within the job
- » Use alternatively the logger from the job in case Talend has created one of instead of create a new one
- » Watching the Log4j configuration file for changes at runtime
- » Access the Log4j configuration via JMX
- » Provide runtime information via JMX
- » Migration to Log4j Version 2